

# E-COMMERCE CLOTHING REVIEW ANALYSIS AND MODEL BUILDING

<sup>1</sup>G. Manikiran, <sup>2</sup>S. Greeshma, <sup>3</sup>P. Vishnu Teja, <sup>4</sup>Y. Sreehari Rao <sup>5</sup>Tanvir H. Sardar

<sup>1,2,3,4</sup>Undergraduate Student, <sup>5</sup>Assistant Professor,

<sup>1,2,3,4,5</sup>Department of CSE (Data Science), Jain Group of Institutions, Bangalore, Karnataka, India

## Abstract:

Understanding customer sentiments is of high importance in marketing strategies today. Not only will it give companies an insight on how customers perceive their products and/or services, but it'll also give them a thought on the way to improve their offers. This paper attempts to know the correlation of various variables in customer reviews on women's clothing e-commerce data, and to classify every review whether it recommends the reviewed product or not and whether it consists of positive or negative or neutral sentiment. To realize these goals, we employed univariate and multivariate analyses on dataset features apart from review texts, which we implemented NLP techniques for recommendation and sentiment classification. Results have shown that a recommendation may be a strong indicator of a positive sentiment score, and vice-versa. On the other hand, ratings in product reviews are fuzzy indicators of sentiment scores. We also found out that the Multinomial Naive Bayes was ready to reach an F1-score of **0.9596** for recommendation classification, and **0.928355** for sentiment classification.

**Keywords:** Data analysis, Machine Learning, NLP, SVM, NLTK, confusion matrix, ROC curve

## Introduction:

With the development of the network, there is an increasing number of people choose to purchase online. Given the shortage of information online, customers are always struggling with issues such as size, quantity, colours, and etc. Therefore, an overview of other customer reviews can help us get a quick impression of the products.

Understanding customer sentiments will enhance the efficiency of online shopping and also give the companies an insight as to how the customers like their product. The best way to improve the customer's experience is by listening to them. The customer's feedbacks always comes in many different forms and languages.

Manually reading all customer's reviews simply wouldn't be possible. The best way to solve this issue is to build a technology model to automatically analysis the customer's feedbacks and the retailer can easily working on improving the customer's experience.

Companies are starting to take social media listening as a tool for understanding their customers, in order to further improve their products and/or services. As a part of this, text analysis has become an active field of research in computational linguistics and natural language processing. One of the most popular problems in this

field is text classification, a task which attempts to categorize documents to one or more classes that may be done manually or computationally.

Towards this direction, In recent years many has shown top interest in classifying sentiments of statements found in social media, review sites, and discussion groups. This task is known as sentiment analysis, a computational technique that uses statistics and natural language processing processes to identify and categorize opinions expressed in a text, particularly, to determine the nature of attitude (positive, negative, or neutral) of the writer towards a topic or a product.

The retailer also can refine sales and marketing strategies or report the important issues that might not be addressed. To achieve these goals, we employed univariate, multivariate analysis and text mining on dataset features and we developed models.

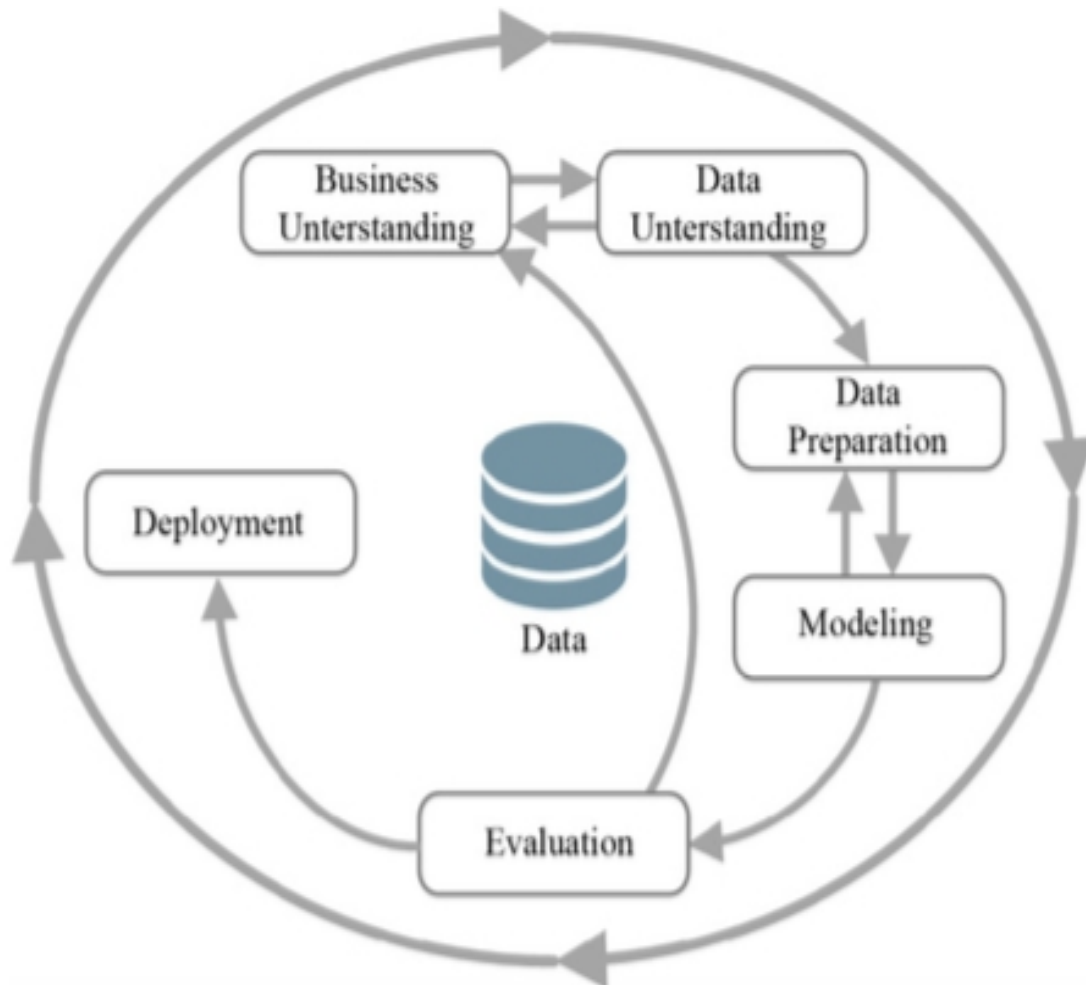
### **Literature Survey:**

[1](Raj Kumar S. Jagdale, Sachin N. Deshmukh Vishal, S. Shirsat; 2019) Dataset for this paper was taken from Amazon in which it has the reviews of laptops, cameras, mobiles, tablets, televisions. After pre-processing is done, to check and classify the reviews into positive or negative they have applied various machine learning algorithms. Finally this paper says that, using machine learning techniques we get best results to classify the Products Reviews. Naïve Bayes got an accuracy of 98.17% whereas for camera reviews support vector machine got an accuracy of 93.54%. [2](Satuluri Vanaja and Meena Belwal; 2018) They used Amazon customer review data in this paper. They have implemented Aspect-level Sentiment Analysis. They have used two machine learning algorithms for classification, SVM classification and NB classification algorithm and the performance is compared based on precision, recall and f1 measure. This paper concludes that we obtained more accuracy from Naïve Bayes algorithm than Support Vector Machine algorithm. [3](P.Sanjay Bhargav, G. Nagarjuna Reddy, R.V. Ravi Chand, K.Pujitha, Anjali Mathur; 2019) In this paper they have used opinion dataset and have implemented various machine learning algorithms using Naïve Bayes Algorithm and opinion Mining Algorithms on the basis of Natural Language Processing. The content-based recommender states the matching of attributes from a particular user profile in which interests and preferences are stored with attributes of content object. If a some morphological variant, is found in the profile and the document, then a match is made and then the document is considered as relevant. [4](KHAN 3et al.; 2019) This paper also proposed a framework which contained data collection, pre-processing, and feature extraction, attribute selection. [5](Sun et al.; 2019) In this paper they proposed a fuzzy product ontology mining algorithm. In this the products are explored from a fine-grained level of online customer reviews. The novel algorithm can not only help a

company to improve their products but also it helps the customers to make better decisions. [6](Yang et al.; 2018)In this paper, in order to figure out the problem with a small number of marker comments, they proposed an evolutionary fuzzy deep belief networks with incremental rules (EFDBNI) algorithm based on fuzzy mathematics and genetic algorithm. The results says that EFDBNI algorithm had a significant improvement over existing methods. This method has also achieved good results in sentiment classification problems with a few labelled comments.

## Materials and Methods

### Architecture:



### Confusion Matrix:

Error matrix which is commonly called a confusion matrix is used for describing the performance of a classification model. The classification model consists of true positives and negatives, false positives, and negatives.

How confused the model is between the classes is shown by the confusion matrix.

	POSITIVE	NEGATIVE	
POSITIVE	True Positive(TN)	False Positive(FN)	Accuracy = (TP + TN)/N Misclassification = (FP + FN)/N Precision = TP/(TP + FP)
NEGATIVE	False Negative(FN)	True Negative(TN)	

FIG-3

### AUC-ROC Curve:

AUC stands for the area under curve and ROC stands for receiver operating characteristics curve.

For checking any classification model's performance AOC-ROC curve is one of the most important evaluation metrics. Multi-class classification problem performance can be checked and visualized using the AUC-ROC curve.

The curve says how much a model is capable of differentiating between classes. The higher the AUC, the better the model is at predicting. When the AUC is higher, better the model is differentiating between Yes and No.

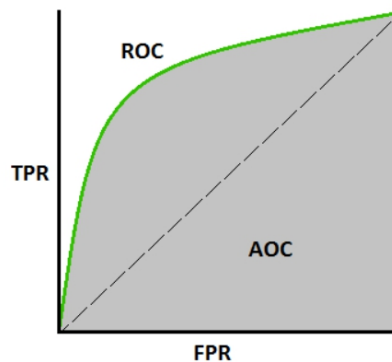


FIG-4

Defining the terms which are used in the above graph:

TPR (True Positive Rate)/Recall/Sensitivity

$$TPR = TP / (TP + FN)$$

$$\text{Specificity} = TN / (TN + FP)$$

$$FPR = 1 - \text{Specificity} = FP / (TN + FP)$$

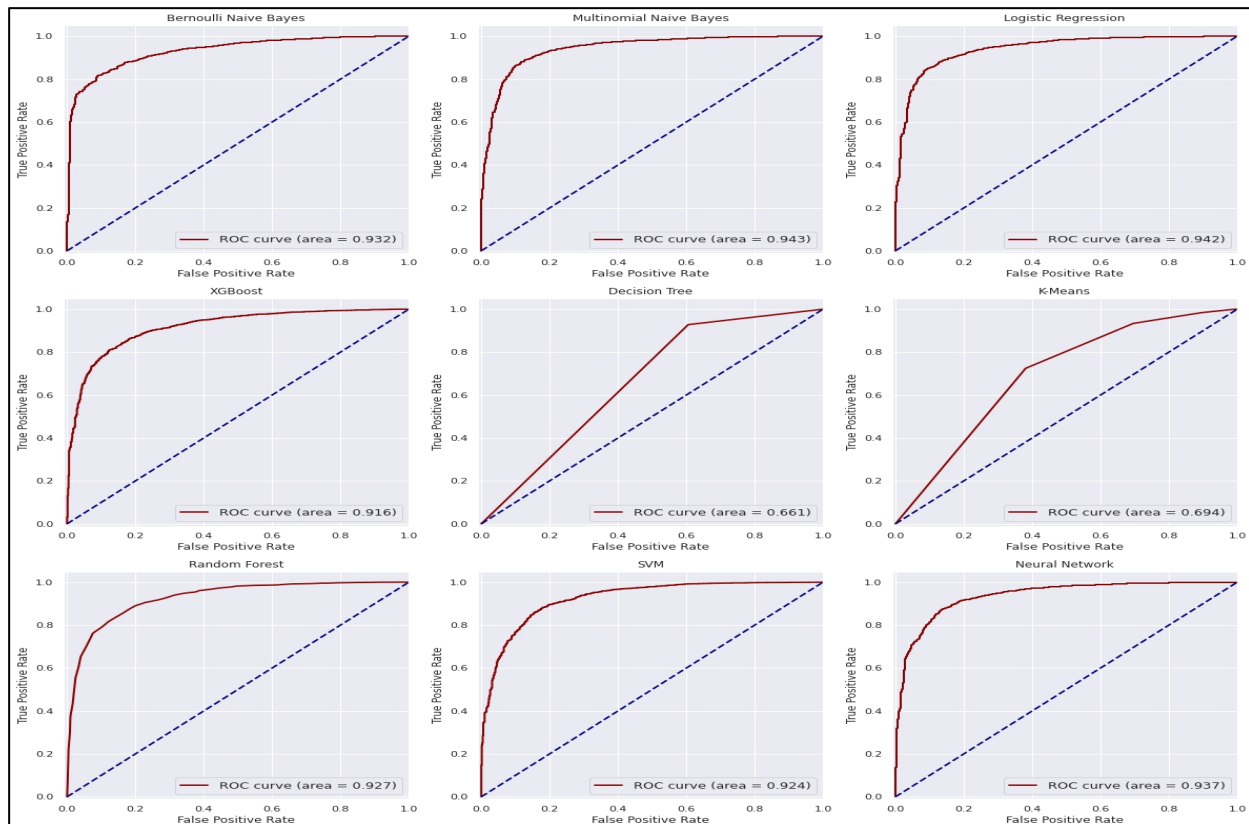
Relationships between sensitivity and specificity are inversely proportional to each other whereas TPR and FPR are directly proportional to each other.

## Result and Conclusion:

### Accuracies and F1 Scores:

Sno.	Algorithm	Accuracy	F1 Score
1	Bernoulli Naive Bayes	0.909182	0.9489
2	Multinomial Naive Bayes	0.928355	0.9596
3	Logistic Regression	0.924318	0.9573
4	XGBoost	0.899764	0.9455
5	Decision Tree	0.862428	0.9220
6	K-Means	0.875714	0.9328
7	Random Forest	0.888664	0.9402
8	SVM	0.920282	0.9555
9	Neural Network	0.922805	0.9562

### ROC and AUC:



In general, an AUC of 0.5 suggests no discrimination, 0.7 to 0.8 is considered acceptable, 0.8 to 0.9 is considered excellent, and more than 0.9 is considered outstanding.

### **Conclusions:**

As the conclusion we have drawn some inferences with the dataset using the visualized charts which may help the industry to increase their production and profit. Below are the inferences with considering the data set as whole:

- 1. Age is not showing much impact in choice selection.*
- 2. General division contributes more reviews.*
- 3. Sleep class has the lowest purchases.*
- 4. Tops department has much following in the current trend.*
- 5. Dresses and knits class tops the list.*
- 6. Dresses class need urgent care.*
- 7. Longue class is most recommended much, it needs to be improved.*
- 8. Good reviews are more (almost 63%).*
- 9. 'Dress' word is the most common word in both positive and negative reviews.*
- 10. Multinomial Naive Bayes model is best suitable for this analysis as it has good performance and accuracy.*

### **Future Scope:**

Analyzing the comments and reviews with the help of python and statistics is a big task for a person who has zero knowledge in those technology, so it's hard to grasp and very difficult to adapt to this product which requires so much of knowledge in python, Visualization etc..., So our future work will be dedicated to the business executives who may use this product as a backend with a supportive easy user interface which will helps them to visualize the data according to their understanding. This may be

customized may be in near future which may depends on the requirement of the customers(Business Firms) and the domain that they are working as the interfaces and attributes differs with respective to the Business place.

If a business executive needs to expose their productivity across the world, he needs to be confident over the product that he is producing, so our product will help them without having any knowledge of python or visualization coding or any tech-expert, he can just use the dataset at the dashboard of the application and draw the charts and inferences with very little efforts. Technology is built by experts and used by people.

Creating a UI and deploying the model will make it easily available for the business executives and thus they can easily access the model and get instant conclusions.

### References:

[1]KHAN, D. M., Rao, T. A. and Shahzad, F. (2019). The classification of customers' sentiment using data mining approaches, *Global Social Sciences Review (GSSR) IV* (IV): 198–212.

[2]Huber, S., Wiemer, H., Schneider, D. and Ihlenfeldt, S. (2019). Dmme: Data mining methodology for engineering applications – a holistic extension to the crisp-dm model, *12th CIRP Conference on Intelligent Computation in Manufacturing Engineering 79*: 403–408.

[3]Pankaj, Pandey, P., Muskan and Soni, N. (2019). Sentiment analysis on customer feedback data: Amazon product reviews, *2019 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon)* pp. 320–322.

[4]Ali, N. M., Hamid, M. M. A. E. and Youssif, A. (2019). Sentiment analysis for movies reviews dataset using deep learning models, *International Journal of Data Mining Knowledge Management Process (IJDKP) 9*(2/3): 19–27.

[5]Jagdale, R. S., Shirsat, V. S. and Deshmukh, S. N. (2019). Sentiment analysis on product reviews using machine learning techniques, *Cognitive Informatics and Soft Computing, Advances in Intelligent Systems and Computing 768* pp. 639–647.

[6]Sun, Q., Niu, J., Yao, Z. and Yan, H. (2019). Exploring ewom in online customer reviews: Sentiment analysis at a fine-grained level, *Engineering Applications of Artificial Intelligence 81*: 68–78.

[7]Yang, P., Wang, D., Du, X.-L. And Wang, M. (2018). Evolutionary dbn for the customers' sentiment classification with incremental rules, *Industrial Conference on*

*Data Mining ICDM 2018: Advances in Data Mining. Applications and Theoretical Aspects* pp. 119– 134.